

PARTIAL TRANSLATION OF JP 58(1983)-193595 A

Publication Date: November 11, 1983

Title of the Invention: TELEPHONE INFORMATION INPUT APPARATUS

Patent Application Number: 57-75282

Filing Date: May 7, 1982

Inventor: K. NAKAMURA

Applicant: HITACHI, LTD.

Claim

1. A telephone information input apparatus characterized in that, in a speech recognition apparatus having a plurality of sets of phoneme standard patterns classified for each speaker and a phoneme sequence word dictionary corresponding to words to be recognized, a first pseudo phoneme pattern for detecting the presence/absence of a push button signal is contained in a particular set of phoneme standard patterns, a second pseudo phoneme pattern for recognizing each push button signal is contained in the remaining sets, and a pseudo phoneme sequence word dictionary corresponding to the first and second pseudo phoneme patterns is provided.

(Page 1, right column, lines 3-10)

According to the push button (PB) signal input 1), among artificial signals obtained by combining speech band sine wave 2 frequencies (one high frequency and one low frequency), a combination of four low frequencies and four high frequencies is currently used with the specification unified, and in principle, 16 kinds of information can be inputted (see Matsuzaka, Uehara, Yatani: A signal system for a push button dial telephone: Nippon Telegraph and Telephone Public Corporation, Telecommunication Research Institute, Research Commercialization Report 17-11, p. 241, November, 1968).

THIS PAGE BLANK (USPTO)

(Page 3, upper right column, line 8 – lower left column, line 12)

The present invention is characterized in that a PB signal is given as one phoneme and word to the speech recognition part 13 shown in FIG. 1 as described above, and the PB signal is detected in the same manner as that of speech recognition, whereby the speech and the PB signal may be used together as information input means. The coexistence of the speech signal and the PB signal is not assumed.

First, simply, the case where frame-based phoneme recognition is performed on a round-robin system with respect to all 16 sets of phoneme patterns will be considered. One PB signal for each of 16 sets (i.e., 16 kinds of PB signals in total) is assigned as a pseudo phoneme standard, and a feature pattern required for the detection thereof may be stored in a phoneme standard pattern memory. The word dictionary corresponding to the PB signals may be configured so as to satisfy the condition that a pseudo phoneme standard with respect to the same PB signal is maintained for a period of time or longer (e.g., 40 milliseconds or longer according the current stipulation) which is required for the reception and detection.

Next, in the case of hierarchical processing in which first-stage recognition is performed by the first two representative clusters, the following is considered.

- 1) It is detected that a signal is a PB signal in the first-stage recognition.

- 2) In the case where it is detected that a signal is a PB signal in the first-stage recognition, it is recognized which PB signal is detected in second-stage recognition.

THIS PAGE BLANK (USPTO)

⑨ 日本国特許庁 (JP)

⑪ 特許出願公開

⑫ 公開特許公報 (A)

昭58—193595

⑤ Int. Cl.³
G 10 L 1/00
H 04 M 1/26
11/00

識別記号

庁内整理番号
7350—5D
7251—5K
7345—5K

④ 公開 昭和58年(1983)11月11日

発明の数 1
審査請求 未請求

(全 5 頁)

⑭ 電話情報入力装置

地株式会社日立製作所中央研究
所内

① 特 願 昭57—75282

⑦ 出 願 人 株式会社日立製作所

② 出 願 昭57(1982)5月7日

東京都千代田区丸の内1丁目5
番1号

③ 発 明 者 中田和男

④ 代 理 人 弁理士 薄田利幸

国分寺市東恋ヶ窪1丁目280番

明 細 書

発明の名称 電話情報入力装置

特許請求の範囲

1. 話者別に分類された複数組の音楽標準ボタンと認識すべき単語に対応した音楽系列単語辞書とを有する音声認識装置において、押しボタン信号の有無を検出するための第1の擬似音楽ボタンを音楽標準ボタンの特定の組の中に有し、残りの組内に個々の押しボタン信号を認識するための第2の擬似音楽ボタンを持ち、第1および第2の擬似音楽ボタンに対応した擬似音楽系列単語辞書を設けたことを特徴とする電話情報入力装置。

発明の詳細な説明

本発明は電話による情報の入力、とくに音声認識を利用した情報入力装置に係り、特にその機能を押ボタン信号による入力の併用にも拡大するの好適な音声認識装置の構成に関する。

従来の電話機による情報の入力手段には次の2つがある。1) 押しボタン信号入力(以下PB入

力と略す)、2) 音声認識入力(以下音声入力と略す)。

1) は音声帯域正弦波2周波(高域、低域各1周波)の組み合わせによる人工的な信号で、現在規格を統一されて使用されているものは低域4周波、高域4周波の組み合わせで原理的に16種類の情報を入力することができる(松坂、上原、矢谷: 押しボタンダイヤル電話用信号方式: 日本電信電話公社電気通信研究所研究実用化報告17-11, P241 昭和43年11月参照)。

この1)の方法によれば情報は確実に入力できるが、情報をすべて数字コードに変換して入力しなければならず、また押しボタン電話機が使えないところでは情報を入力することができない。

2) は音声認識によつて、音声のままて情報を直接入力しようとするもので、便利ではあるが、常に確実、正確に情報が入力できるとは限らない(長島、中津: 音韻単位の標準ボタンを用いた長時間単語音声認識装置、日本音響学会音声研究会資料、878-22, 1979, 渡辺、直理、千葉

他；不特定話者用音声認識装置SR-1000シリーズ，日本音響学会講演論文集，3-1-24，1981年5月）

本発明の目的は、従来の音声認識装置の構成を基本とし、これにごくわずかの追加を行うことによつて、あらかじめ音声信号かPB信号かがわからなくても、それぞれ認識が行なわれ、そのことによつて音声とPB信号を自由に併用して使用でき、電話機による情報入力機能を拡大する手段を提供することにある。

まず、従来の電話情報入力用音声認識システムの構成を第1図に示す。

第1図において、加入者電話機11から交換機12を通つた音声信号121は音声認識部13に入力され、業務処理部14からの認識要求信号141を受けてその認識処理をおこなう。主業務処理部14では、認識結果を確認するために認識完了信号142を受けて音声出力部15に出力要求信号151を送出し、音声出力の終了を出力完了信号152により確認する。

と単語辞書メモリ24中に格納されている標準単語（たとえば、単語番号1，2，……に対応してそれぞれ音楽記号系列i・chi，ni，……などで表わされる単語）との非線形マッチング演算がDPマッチング部25においておこなわれ、その結果得られた距離和251の大きさにもとづいて単語判定部26で入力音声の判定がおこなわれ、認識結果27が出力される。

この認識処理の特徴は、電話入力された不特定話者の音声認識を、16組の音楽標準パターンによるフレーム別認識をおこなう第1段と、フレーム別認識の結果と音楽記号系列単語辞書とのDPマッチングをおこなう第2段とからなる2段のボタン整合に分解し、第1段では音楽標準パターンにたいして話者の音声波形における音響的な特性にもとづいて16組のクラスタリング（組み合わせ）をおこない、第2段では1つの単語に対して複数の音楽記号系列単語辞書をもうけて、発話の変化、たとえば母音の無声化や鼻音化、にに対処していることである。本方式はこの2段処理によつて

一万、交換機12からの応答信号122を受けて発信制御部16から出力された応答信号161が主業務処理部14に入力されると、電文処理部17にたいして送信要求信号143が送出される。

これを受けた電文処理部17は通信制御部18にたいし、送信要求信号171を送ることによりリレーコンピュータ19から発せられ通信制御部18を通つた電文181を受信して発信制御部16にたいし発信要求信号172を送り信号162を発信させる。

第2図は第1図における音声認識部13のブロック構成を示す。

第2図(b)で示す波形の入力音声20（ichi）から音声分析部21において抽出された特徴パラメータの系列211と音楽標準パターンメモリ22中に格納されている例えば16組の音楽（a，i，……など）の特徴パラメータ（最尤スペクトルパラメータ、LPCケプストラム係数など）との距離が距離計算部23において計算される。

距離計算部23から出力された距離の系列231

所要メモリ量が少なくかつ不特定話者の音声に対して高い認識能力を持つことが知られている。

「一桁の数字音（0～9の10語）」および、「はい」、「いいえ」、「どうぞ」、「もう一度」、「ほりゆう（保留）」、「とりけし（取り消し）」の6語を含む16語に対して620名の男女による認識結果の一例を表1に示す（電電公社通信研究所発表）。

表1 男女別の誤り率[%]

	男	女	平均
尤度による距離	4.45	3.22	3.86
LPCケプストラム距離	4.57	3.18	3.90

なお、表1で距離尺度としてとられているのは、音声認識のための特徴として使われるパラメータの一例であり、このいづれを用いても誤り率はほとんど違わないことをあらわしている。

この方式のもう一つの特徴は、16組に分類された多数の（最大40個程度）音楽標準パターンとの整合によつて、フレーム別に音楽系列を認識

し、その結果と単語音素系列との比較によつて単語を認識するに当つて、その処理量を軽減し、実時間認識を可能にするため、そのフレーム別音素認識を第3図(a)に示すように2段に分けて階層的に行つてゐることである。すなわち、16組のボタンの中、男声の代表として作られている例えば第1の組と、女声の代表として作られている例えば第16の組との2組の標準パターンで、まず第1段の認識を行い、その中で整合度の良いものN語をえらび、そのN語に対象を限定し、改めて上記16組の音素標準パターンのすべてを使つて再認識を行う。Nの値としては第3図(b)に示す実験結果から $N=4$ にとれば、誤認識による誤りが少なく、処理量(計算量)も少なくて済むことがわかる。ここで計算量の比率とは、

$$\frac{\text{一次選択ありの計算量}}{\text{一次選択なしの計算量}} \times 100 \text{ で与えられる。}$$

結果的に16語に対して16語 \times 16語=256組 \times 語の処理を、2組 \times 16語+16組 \times 4語=96組 \times 語の処理に軽減している。

応する単語辞書としては、受信検出しなければならないとされている時間以上(たとえば現行規定によれば40ミリ秒以上)同一のPB信号に対する擬音素標準が維持するという条件を満足するように構成すればよい。

次に最初2個の代表クラスターによつて第1段目の認識が行なわれるという階層処理の場合には次のように考える。

1) 第1段目の認識でPB信号であることを検出する。

2) 第1段目認識でPB信号と検出された場合第2段目でそのいずれであるかを認識する。

以下さらに具体的に説明する。

音韻認識において、LPC(線形予測)分析にもとづいて尤度比による整合をとる場合について考える。

0.3 kHzから3.4 kHzまでに帯域制限された音声信号に対して、通常 $p=10$ 次の分析が行なわれる。

この分析の結果、原理的には $p/2$ 個のスペク

さて通常のPB信号は、いわゆるPB信号受信器で検出される。通常の使い方では情報を入力する信号の形式がPB信号であるか音声信号であるかはあらかじめ決まつており、分離して行なわれる。例えば通常の使い方ではPB信号は情報センターへのアプローチに使われ、第1図における発信制御部13で受信検出される。

本発明のポイントはすでに述べた第1図の音声認識部13へ、音素および単語の1個としてPB信号を加え、音声認識と全く同じ形式でPB信号を検出することによつて、音声とPB信号を情報入力手段として併用してもよいようにしようとするものである。ただし音声信号とPB信号の同時共存は仮定しない。

まず簡単に、16組の音素パターンのすべてと総当たりでフレーム別音素認識が行なわれる場合を考える。このときは16組の各組に1個、あわせて16種のPB信号を擬音素標準として割り当て、その検出に必要な特徴パターンを音素標準パターンメモリに記憶させておけばよい。PB信号に対

トルの共振周波数いわゆるホルマント周波数が指定される。すなわち $p=10$ の場合、5個の周波数を指定することができる。この5個の周波数を、低、高の両周波数帯に、第4図に示すように割り当てれば、16個の周波数の中の任意の6個をカバーするように設定することができ、2組によつて任意の12個をカバーするようにすることができる。

第4図において、1, 2, 3, 4, 5は $\Phi 1$ クラスターの割り当て周波数を示し、I, II, III, N, Vは $\Phi 1b$ クラスターの割り当て周波数を示す。

PB信号としては、16個の中から実際には10数字と制御用に2個(たとえば・印と Φ 印)が用いられるのが普通であり、12個を検出できればよい。日本国内では低域は4周波(697, 770, 852, 941 Hz)であるが、高域は3周波(1209, 1336, 1497 Hz)しか使っていない。

これらを検出するパラメータは次式から導出することができる。

指定周波数を $\{f_i\} = (f_1, f_2, f_3, \dots, f_4, f_5)$ とするとき

$$\begin{cases} \beta_i = r_i e^{j\theta_i} & r_i = e^{-\alpha_i T} \\ \bar{\beta}_i = r_i e^{-j\theta_i} & \theta_i = 2\pi f_i T \end{cases}$$

ここで T はサンプリング周期、 b_i は f_i の共振帯域幅であり、PB信号の場合、許容信号周波数変動幅は $\pm 2\%$ と規定されているから $b_i = f_i \times 4\%$ 程度にとればよい。

これから

$$\begin{aligned} & (Z - \beta_1)(Z - \bar{\beta}_1)(Z - \beta_2)(Z - \bar{\beta}_2) \\ & (Z - \beta_3)(Z - \bar{\beta}_3) \dots \dots \dots (1) \end{aligned}$$

の10次方程式を作り、それを

$$Z^{10} + \alpha_1 Z^9 + \alpha_2 Z^8 + \dots + \alpha_9 Z + \alpha_{10} = 0 \quad (2)$$

において(1)式と(2)式の Z の等べき係数を $\alpha_1, \dots, \alpha_{10}$ とおけば、 $(\alpha_1, \dots, \alpha_{10})$ が求められる。

音韻標準パターンとして使われる逆スペクトル係

数は、この α の系列に α_0 とした1を加えた系列の相関係数として、

$$\begin{aligned} A_0 &= 1 + \alpha_1 + \alpha_2 + \dots + \alpha_{10} \\ A_1 &= \alpha_1 + \alpha_1 \alpha_2 + \alpha_2 \alpha_3 + \dots + \alpha_9 \alpha_{10} \\ &\vdots \\ A_9 &= \alpha_9 + \alpha_9 \alpha_{10} \\ A_{10} &= \alpha_{10} \end{aligned}$$

と求められる。

※2から※15までのクラスに記憶される音素パラメータは、現実の個々のPB信号をLPC分析することによつて求めることができる。

なお実際には12個のPB信号をすべて対象とする必要はなく、第4図にその1例を示すように、※1クラスによつて6種類、※16クラスによつて6種類が指定されるから、この6種類についてのみ第2段の認識実験を行えばよい。

本発明の一実施例を第2図を用いて説明する。

入力音声20(擬似音声波形としてPB信号であることもある)は音声分析部21で相関係数 $\{r_i^{(n)}\}$ の算出とLPC(線形予測)分析がさ

れ、残差電力 E_n が計算される。

次に距離計算部23で各フレーム毎に音素標準パターン $\{A_i^{(n)}\}$, $i = 0 \sim 10$, $n = 1 \sim 8$ と入力 X の相関係数 $\{r_i^{(n)}\}$, $i = 0 \sim 10$ と E_n によつて次式によつて尤度比が計算される。

$$L_n = [A_0 \cdot r_0 + 2 \sum_{i=1}^{10} A_i \cdot r_i] / E_n \quad (3)$$

この L_n を尺度とする入力音素系列マトリックスと音素記号系列単語辞書との間でDPマッチングによる整合がとられ、最適整合のものが認識結果として出力される。その場合、すでに説明したように16組の音韻クラスにおいて、代表的な2つ、たとえば※1(男声代表)と※16(女声代表)のみを用いた第1段の認識が行われ、候補単語が N 個にしばらくされる。このとき、※1と※16のクラスターに追加されたPB信号検出用のパターンによつて第1候補がPB信号であると検出されたときは、 N 個の候補として、12種類のPB信号の中の6個を候補として第2段目の認識

を行う。その他は従来の音声認識と全く同じである。

この場合、個別PB信号に対応する擬似音素パターンを2組から15組に1個ずつ加えないで、第17組としてPB信号用のクラスを構成すれば、第一段目でPB信号として検出されたときは、このクラスについてのみフレーム別音素認識を行えばよいように構成することもできる。

また一般に行なわれている単語レベルでの複数標準パターンによる音声認識においては、PB信号に対しては16組クラスター総当りで説明した一段目の認識のみでよいことは自明である。

以上説明したように本発明によれば、音声とPB信号を何ら区別することなく電話による情報入力手段として利用することができ、音声入力の簡便さとPB入力の確実さの特色を活かした情報入力が可能となる。

たとえば、音声によつては比較的長く、文脈効果の利用しやすい制御語のみを入力し、短かくて文脈効果の利用しえない数字データはPB入力と

するといった使い方も可能となる。

あるいはPB電話機を利用できる人には確実なPB入力を、利用できない人には音声入力を使うシステムをサービスすることもできる。

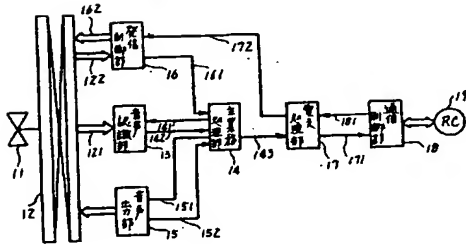
図面の簡単な説明

第1図は従来音声認識応答システムの構成図、第2図はその音声認識部の説明図、第3図は実際に行なわれている階層認識処理の説明図、第4図はPB信号検出用擬似音韻パターンによる検出可能領域の説明図である。

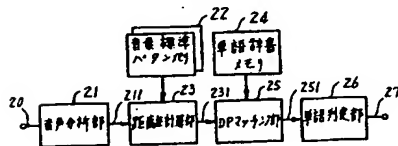
1 3…音声認識部。

代理人 弁理士 薄田利幸

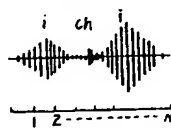
第 1 図



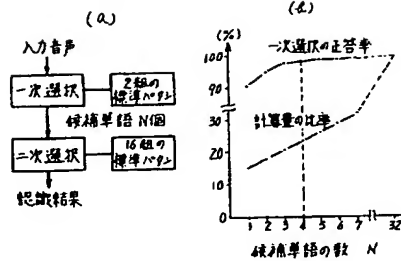
第 2 図
(A)



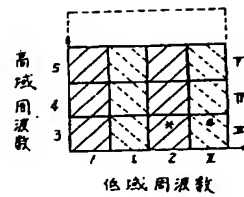
(B)



第 3 図



第 4 図



THIS PAGE BLANK (USPTO)